

DimmWitted: A Study of Main-Memory Statistical Analytics

Ce Zhang^{†‡} Christopher Ré[‡]

[†]University of Wisconsin-Madison

[‡]Stanford University

{czhang, chrismre}@cs.stanford.edu

ABSTRACT

We perform the first study of the tradeoff space of access methods and replication for statistical analytics executed in the main memory of a Non-Uniform Memory Access (NUMA) machine. Our goal is to understand tradeoffs in accessing the data in row- or column-order and at what granularity one should share the model and data for a statistical task. Statistical analytics systems differ from conventional SQL-analytics in the amount and types of memory incoherence they can tolerate. We study this new tradeoff space, and discover there are tradeoffs between hardware and statistical efficiency. We argue that our tradeoff study may provide valuable information for designers of analytics engines: for each system we consider, our prototype engine can run at least one popular task at least 100× faster. We conduct our study across five architectures using popular models including support vector machines, logistic regression, Gibbs sampling, and neural networks.

1. INTRODUCTION

Statistical data analytics is one of the hottest topics in data-management research and practice. Today, even small organizations have access to machines with large main memories (via Amazon’s EC2) or for purchase at \$5/GB. As a result, there has been a flurry of activity to support main-memory analytics in both industry (Google Brain, Impala, and Pivotal) and research (GraphLab, and MLlib). Each of these systems picks one design point in a larger tradeoff space. The goal of this paper is to define and explore this space. We find that today’s research and industrial systems under-utilize commodity modern hardware for analytics—sometimes by two orders of magnitude. We hope that our study identifies some useful design points for the next generation of such main-memory analytics systems.

Our study examines analytics on commodity multi-socket, multi-CPU, non-uniform memory access (NUMA) machines, which are the de facto standard machine configuration and so a natural target for an in-depth study. Moreover, our

experience with several enterprise companies suggests that, after appropriate preprocessing, a large class of enterprise analytics problems fit into the main memory of a single, modern machine. While this architecture has been recently studied for traditional SQL-analytics systems [5], it has not been studied for *statistical* analytics systems. We believe that our study is the first such study.

Statistical analytics systems are different from traditional SQL-analytics systems. In comparison to traditional SQL-analytics, the underlying methods are intrinsically robust to error. On the other hand, traditional statistical theory does not consider which operations can be efficiently executed. This leads to a fundamental tradeoff between *statistical efficiency* (how many steps are needed until convergence to a given tolerance) and *hardware efficiency* (how efficiently those steps can be carried out). This paper studies this tradeoff space.

To describe such tradeoffs more precisely, we describe the typical setup of a typical analytics task. The input data is a matrix in $\mathbb{R}^{N \times d}$ and the goal is to find a vector $x \in \mathbb{R}^d$ that minimizes some (convex) loss function, say the logistic loss or hinge loss (SVM). Typically, one makes several complete passes over the data while updating the model; we call each such pass an *epoch*. There may be some communication at the end of the epoch, e.g., in bulk-synchronous parallel systems like Spark. We identify three tradeoffs that have not been explored in the literature: (1) *access methods for the data*, (2) *model replication*, and (3) *data replication*. Current systems have picked one point in this space;¹ we explain each space and discover points that have not been previously considered. Using these new points, we find that we can perform 100× faster than previously explored points in the tradeoff space for several popular tasks.

Access Methods. Analytics systems access (and store) data in either row-major or column-major order. For example, systems that use *stochastic gradient methods* (SGD) access the data row-wise; examples include MADlib [8] in Impala and Pivotal, Google Brain [13], and MLlib in Spark [19]; and *stochastic coordinate descent* (SCD) access the data column-wise; examples include GraphLab [17], Shogun [25], and Thetis [27]. These methods have essentially identical statistical efficiency (see Section 3.2), but their wall-clock performance can be radically different due to hardware efficiency. However, this tradeoff has not been systematically studied. To study this tradeoff, we introduce a storage abstraction

¹We participated in MADlib over Pivotal Greenplum [8] and Impala [10], and a product from Oracle [9].

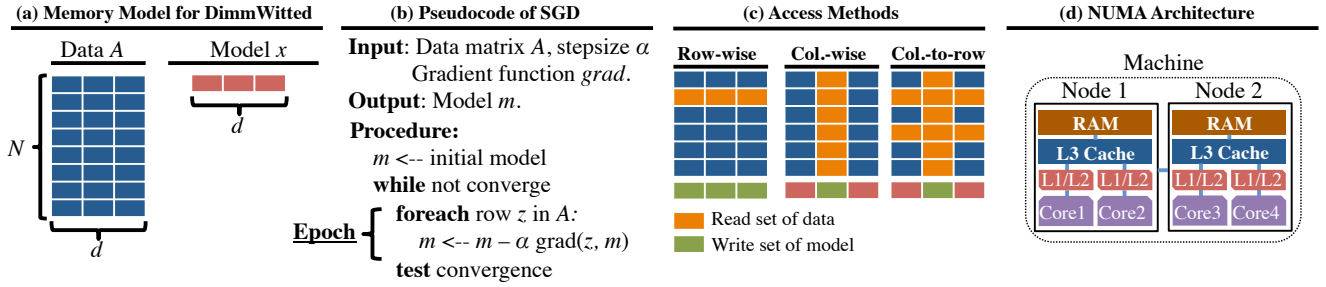


Figure 1: Illustration of DimmWitted’s Memory Model, Pseudocode for SGD, Different Statistical Methods in DimmWitted and Their Access Patterns, and NUMA Architecture.

that captures the access patterns of popular statistical analytics tasks and a prototype called DIMMWITTED. In particular, we identify three access methods that are used in popular analytics tasks including standard supervised machine learning models like SVMs, logistic regression, least squares; and more advanced methods like neural networks and Gibbs sampling on factor graphs. For different access methods for the same problem, we find that the time to converge to a given loss can differ by up to $100\times$.

We also find that no access method dominates all others, and thus an engine may be designed to include both access methods. To show that it may be possible to support both methods in a single engine, we develop a simple cost model to choose among these access methods. We describe a simple cost model that selects a nearly optimal point in our data sets, models, and different machine configurations.

Data and Model Replication. We study two sets of trade-offs: the level of granularity, and the mechanism by which mutable state and immutable data are shared in analytics tasks. We describe the tradeoffs we explore in both (1) mutable state sharing, which we informally call *model replication*, and (2) data replication.

(1) Model Replication. During execution, there is some state that the task mutates (typically an update to the model). We call this state, which may be shared among one or more processors, a *model replica*. We consider three different granularities at which to share model replicas:

- The **PerCPU** approach treats a NUMA machine as a distributed system in which every CPU is treated as an individual machine, e.g., in bulk-synchronous models like MLlib on Spark or event-driven systems like GraphLab. These approaches are the classical shared-nothing and event-driven architectures, respectively. In PerCPU, the part of the model that is updated by each CPU is only visible to that CPU until the end of an epoch. This method is efficient and scalable from a hardware perspective, but it is less statistically efficient as there is only coarse-grained communication between CPUs.
- The **PerMachine** approach treats each processor as if it has uniform access to memory. This approach is taken in Hogwild! and Google Downpour [6]. In this method, the hardware takes care of coherence of the shared state. The PerMachine method is statistically efficient due to high communication rates, but it may

cause contention in the hardware, which may lead to suboptimal running times.

- A natural hybrid is **PerNode**; this method uses the fact that PerCPU communication through the last-level cache (LLC) is dramatically faster than communication through remote main memory. This method is a novel hybrid of statistical and hardware efficiency; for some models, PerNode can be an order of magnitude faster.

Because model replicas are mutable, a key question is: *how often should we synchronize model replicas?* We find that it is beneficial to synchronize the models as much as possible—so long as we do not damage throughput to data in main memory. A natural idea, then, is to use **PerMachine** sharing, in which the hardware is responsible for synchronizing the replicas. However, this decision can be suboptimal as the cache-coherence protocol may stall a processor to preserve coherence—but this information may not be worth the cost of a stall from a statistical efficiency perspective. We find that the **PerNode** method, coupled with a simple technique to batch writes across sockets, can dramatically reduce communication and CPU stalls. The **PerNode** method can result in an over $10\times$ runtime improvement. This technique depends on the fact that we do not need to maintain the model consistently: we are effectively delaying some updates to reduce the total number of updates across sockets (which lead to processor stalls).

(2) Data Replication. The data for analytics is immutable, so there are no synchronization issues for data replication. The classical approach is to partition the data to take advantage of higher aggregate memory bandwidth. However, each partition may contain skewed data, which may slow convergence. Thus, an alternate approach is to fully replicate the data (say, per NUMA node). In this approach, each node accesses that node’s data in a different order, which means that the replicas provide non-redundant statistical information; in turn, this reduces the variance of the estimates based on the data in each replicate. We find that on some tasks, fully replicating the data four ways can be almost $4\times$ faster than the sharding strategy to converge to the same loss. The full replication approach would be nonsensical in SQL-analytics, as it would be pure overhead.

Summary of Contributions. We are the first to study the three tradeoffs above for main-memory statistical analytics

systems. An outcome of our study is that treating NUMA-machines as distributed systems or shared-memory systems is suboptimal; this includes many current research and industrial approaches. Fundamentally, computing systems are becoming increasingly heterogeneous due to the end of voltage scaling, and this is a first study of how to deal with such asymmetry for analytics systems. We design a storage manager, DIMMWITTED, that shows it is possible to exploit these ideas on real data sets. Finally, we evaluate our techniques on multiple real datasets, models, and architectures.

2. BACKGROUND

In this section, we describe the memory model for DIMMWITTED, which provides a unified memory model to implement popular analytics methods. Then, we recall some basic properties of modern NUMA architectures.

Data for Analytics. The data for an analytics task is a pair (A, x) , which we call the data and the model respectively. For concreteness, we consider a matrix $A \in \mathbb{R}^{N \times d}$. In machine learning parlance, each row is called an *example*. Thus, N is often the number of examples and d is often called the dimension of the model. There is also a model, typically a vector in $x \in \mathbb{R}^d$. The distinction is that the data A is read-only while the model vector, x , will be updated during execution. From the perspective of this paper, the important distinction we make is that data is an immutable matrix while the model (or portions of it) are mutable data.

Iterative Analytics Algorithms. All algorithms we consider make several passes over the data; we refer to each such pass as an *epoch*. A popular example algorithm is Stochastic Gradient Descent (SGD), which is widely used in web-companies, e.g., Google Brain [13] and VowPal Wabbit [1], and in enterprise systems like Pivotal, Oracle, and Impala. Pseudocode for this method is shown in Figure 1(b). During each epoch, SGD reads a single example z ; it uses the current value of the model and z to estimate the derivative; it then updates the model vector with this estimate. It reads each example in this loop. After each epoch, these methods test convergence (usually by computing or estimating the norm of the gradient); this computation requires a scan over the complete dataset.

2.1 Memory Models for Analytics

We design DIMMWITTED’s memory model to capture the trend in recent high performance sampling and statistical methods. There are two aspects to this memory model: the *coherence level* and the *storage layout*.

Coherence Level. Classically, memory systems are coherent: reads and writes are executed atomically. For analytics systems, we say a memory model is *coherent* if reads and writes of the entire model vector are atomic. That is, access to the model is enforced by a critical section. However, many modern analytics algorithms are designed for an *incoherent* memory model. The Hogwild! method showed that one can run such a method in parallel without locking but still provably converge. The Hogwild! memory model relies on the fact that writes of individual components are atomic, but it does not require that the entire vector be updated atomically. However, atomicity at the level of the cacheline is

Algorithm	Access Method	Implementation
Stochastic Gradient Descent	Row-wise	MADlib, Spark, Hogwild!
Stochastic Coordinate Descent	Column-wise	GraphLab, Shogun, Thetis
	Column-to-row	

Figure 2: Algorithms and Their Access Methods.

provided by essentially all modern processors. Empirically, these results allow one to forgo costly locking (and coherence) protocols. Similar algorithms have been proposed for other popular methods including Gibbs sampling [11, 24], Stochastic Coordinate Descent (SCD) [22, 25], and linear systems solvers [27]. This technique has been applied by Dean et al. [6] to solve convex optimization problems with billions of elements in a model. This memory model is distinct from the classical, *fully coherent* database execution.

The DIMMWITTED prototype allows us to specify that a region of memory is coherent or not. This region of memory may be shared by one or more processors. If the memory is only shared per thread, then we can simulate a shared-nothing execution. If the memory is shared per machine, we can simulate Hogwild!.

Access Methods. We identify three distinct access paths used by modern analytics systems, which we call row-wise, column-wise, and column-to-row. They are graphically illustrated in Figure 1(c). Our prototype supports all three access methods. All of our methods perform several epochs, that is, passes over the data. However, the algorithm may iterate over the data row-wise or column-wise.

- In *row-wise access*, the system scans each row of the table and applies a function that takes that row, applies a function to it, and then updates the model. This method may write to all components of the model. Popular methods that use this access method include stochastic gradient descent, gradient descent, and higher order methods (like l-BFGS).
- In *column-wise access*, the system scans each column j of the table. This method reads just the j component of the model. The write set of the method is typically a single component of the model. This method is used by stochastic coordinate descent.
- In *column-to-row access*, the system iterates conceptually over the columns. This method is typically applied to sparse matrices. When iterating on column j it will read all rows in which column j is non-zero. This method also updates a single component of the model. This method is used by non-linear support vector machines in GraphLab and is the de facto approach for Gibbs sampling.

The system is free to iterate over rows or columns in essentially any order (although typically some randomness in the ordering is desired). Figure 2 classifies popular implementations by their access method.

2.2 Architecture of NUMA Machines

Name (abbrv.)	#Node	#Cores/Node	RAM/Node (GB)	CPU Clock (GHz)	LLC (MB)
local2 (l2)	2	6	32	2.6	12
local4 (l4)	4	10	64	2.0	24
local8 (l8)	8	8	128	2.6	24
ec2.1 (e1)	2	8	122	2.6	20
ec2.2 (e2)	2	8	30	2.6	20

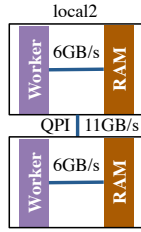


Figure 3: Summary of Machines and Memory Bandwidth on local2 Tested with STREAM [4].

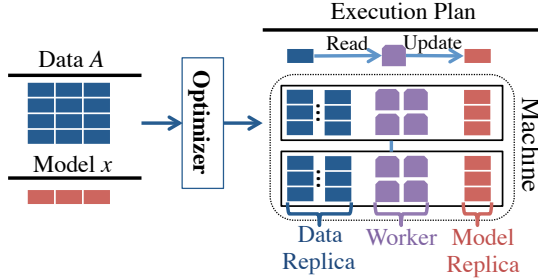


Figure 4: Illustration of DimmWitted's Engine.

We briefly describe the architecture of a modern NUMA machine. As illustrated in Figure 1(d), a NUMA machine contains multiple NUMA nodes. Each node has multiple CPUs (cores) and processor caches, including the L3 cache. Each node is directly connected to a region of DRAM. NUMA nodes are connected to each other by buses on the main board; in our case, this connection is the Intel Quick Path Interconnects (QPIs), which has a bandwidth as high as 25.6GB/s.² To access DRAM regions of other NUMA nodes, data is transferred across NUMA nodes using the QPI. These NUMA architectures are cache coherent, and the coherency actions use the QPI. Figure 3 describes the configuration of each machine that we use in this paper. Machine controlled by us have names with the prefix ‘local’; the other machines are Amazon EC2 configurations.

3. THE DIMMWITTED ENGINE

We describe the tradeoff space that DIMMWITTED’s optimizer considers, namely (1) Access Method Selection, (2) Model Replication, and (3) Data Replication. To help understand the statistical-versus-hardware tradeoff space, we present some experimental results in a *Tradeoffs* paragraph within each subsection. We also describe key implementation details for DIMMWITTED.

3.1 System Overview

For each analytics task that we study, we assume the user provides data $A \in \mathbb{R}^{N \times d}$ and an initial model that is a vector of length \mathbb{R}^d . In addition, for each access method listed above, there is a function of an appropriate type that solves the same underlying model. For example, we provide

²www.intel.com/content/www/us/en/io/quickpath-technology/quick-path-interconnect-introduction-paper.html

Tradeoff	Strategies	Existing Systems
Access Methods	Row-wise	SP, HW
	Column-wise	GL
	Column-to-row	GL, SP
Model Replication	Per CPU	GL, SP
	Per Node	HW
	Per Machine	GL, SP, HW
Data Replication	Sharding	
	Full Replication	

Figure 5: A Summary of DimmWitted’s Tradeoffs and Existing Systems (GraphLab (GL), Hogwild! (HW), Spark (SP)).

both a row- and column-wise way of solving a support vector machine. Each method takes two arguments, with the first acting as a pointer to a model (we use this indirection below).

- f_{row} captures the the row-wise access method, and its second argument is the index of a single row.
- f_{col} captures the column-wise access method, and its second argument is the index of a single column.
- f_{ctr} captures the column-to-row access method, and its second argument is a pair of one column index and a set of row indexes. These rows correspond to the non-zero entries in a data matrix for a single column.³

Each of the functions modify the model to which they receive a pointer in place. However, in our study f_{row} can modify the whole model, while f_{col} and f_{ctr} only modify a single variable of the model. We call the above tuple of functions a *model specification*. Note that a model specification contains either f_{col} or f_{ctr} but typically not both.

Execution. Given a model specification, our goal is to generate an execution plan. An execution plan, schematically illustrated in Figure 4, specifies three things for each CPU core in the machine: (1) a subset of the data matrix to operate on, (2) a replica of the model to update, and (3) the access method used to update the model. We call the set of replicas of data and models *locality groups* as the replicas are described physically, i.e., they correspond to regions of memory that are local to particular NUMA nodes, and one or more workers may be mapped to each locality group. The data assigned to distinct locality groups may overlap. We use DIMMWITTED’s engine to explore three tradeoffs:

- (1) **Access Methods** in which we can select between either the row or column method to access the data.
- (2) **Model Replication** in which we choose how to create and assign replicas of the model to each worker. When a worker needs to read or write the model, it will read or write the model replica that it is assigned.
- (3) **Data Replication** in which we choose a subset of data tuples for each worker. The replicas may be overlapping, disjoint, or some combination.

³Define $S(j) = \{i : a_{ij} \neq 0\}$. For a column j , the input to f_{ctr} is a pair $(j, S(j))$.

Algorithm	Read	Write (Dense)	Write (Sparse)
Row-wise	$\sum n_i$	dN	$\sum n_i$
Column-wise	$\sum n_i$	d	
Column-to-row	$\sum n_i^2$		

Figure 6: Per Epoch Execution Cost of Row- and Column-wise Access. Given a data set $A \in \mathbb{R}^{N \times d}$, let n_i be the number of non-zero elements a_i .

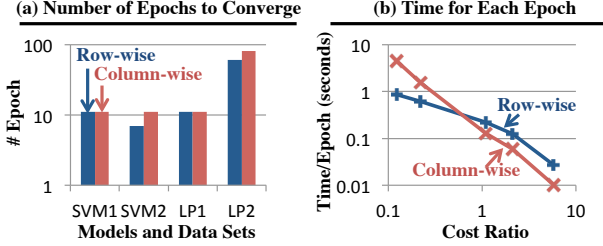


Figure 7: Illustration of the Method Selection Tradeoff. (a) These four data sets are RCV1, Reuters, Amazon, and Google, respectively. (b) The “Cost Ratio” is defined as the ratio of costs estimated for row-wise and column-wise methods: $(1 + \alpha) \sum_i n_i / (\sum_i n_i^2 + \alpha d)$, where n_i is the number of non-zero elements of i^{th} row of A and α is the cost ratio between writing and reading. We set $\alpha = 10$ to plot this graph.

Figure 5 summarizes the tradeoff space and how DIMMWITTED chooses among the strategies that we propose. In each section, we illustrate the tradeoff along two axes, namely (1) the *statistical efficiency*, i.e., the number of epochs it takes to converge; and (2) *hardware efficiency*, the time that each method takes to finish a single epoch.

3.2 Access Method Selection

In this section, we examine each different access method: row-wise, column-wise, and column-to-row. We find that the execution time of an access method depends more on hardware efficiency than statistical efficiency.

Tradeoffs. We consider the two tradeoffs that we use for a simple cost model (Figure 6). Let n_i be the number of non-zeros in row i , so that the number of reads in a row-wise access method is $\sum_{i=1}^N n_i$. Since each example is likely to be written back in a dense write, we perform dN writes per epoch. Our cost model combines these two costs linearly with a factor α that accounts for writes being more expensive on average, due to contention. The factor α is estimated at installation time by measuring on a small set of data sets. The parameter α is in 4 to 12 and grows with the number of sockets, e.g., for local2, $\alpha \approx 4$, and local8, $\alpha \approx 12$. Thus, α may increase in the future.

Statistical Efficiency. We observe that each access method has comparable *statistical efficiency*. To illustrate this, we run all methods on all of our data sets and report the number of epochs that one method converges to a given error to the optimal loss, and Figure 7(a) shows the result on four data sets with 10% error. We see that the gap of

number of epochs cross different methods are small (always within 50% of each other).

Hardware Efficiency. Different access methods can change the time per epoch by up to a factor of 10 \times , and there is a cross-over point. To see this, we run both methods on a series of synthetic data sets where we control the number of non-zero elements per row by subsampling each row on the Music data set (see Section 4 for more details). For each subsampled data set, we plot the cost ratio on the x -axis, and we plot their actual running time per epoch in Figure 7(b). We see a cross-over point on the time used per epoch: when the cost ratio is small, row-wise outperforms column-wise by 6 \times as the column-wise method reads more data; on the other hand, when the ratio is large, the column-wise method outperforms the row-wise method by 3 \times as the column-wise method has lower write contention. We observe similar cross-over points on our other data sets.

Cost-based Optimizer. DIMMWITTED estimates the execution time of different access methods using the number of bytes that each methods reads and writes in one epoch, as shown in Figure 6. For writes, it is slightly more complex: for models like SVM, each gradient step in row-wise access only updates the coordinates where the input vector contains non-zero elements. We call this scenario a *sparse* update; otherwise it is a *dense* update.

DIMMWITTED needs to estimate the ratio of the cost of reads to writes. To do this, it runs a simple benchmark dataset. We find that on all of our eight datasets, five statistical models, and five machines that we used in the experiments, the cost model is robust to this parameter: as long as writes are 4 \times to 100 \times more expensive than reading, the cost model makes the correct decision between row-wise and column-wise access.

3.3 Model Replication

In DIMMWITTED, we consider three model replication strategies. The first two strategies, namely PerCPU and PerMachine, are similar to traditional shared-nothing and shared-memory architecture, respectively. We also consider a hybrid strategy, PerNode designed for NUMA machines.

3.3.1 Granularity of Model Replication

The difference between the three model replication strategies is the granularity of replicating a model. We first describe PerCPU and PerMachine and their relationship with other existing systems (Figure 5). We then describe PerNode, a simple, novel hybrid strategy that we designed to leverage the structure of NUMA machines.

PerCPU. In the PerCPU strategy, each CPU maintains a mutable state, and these states are combined to form a new version of the model (typically at the end of each epoch). This is essentially a shared-nothing architecture; it is implemented in Impala, Pivotal, and Hadoop-based frameworks. PerCPU is popularly implemented by state-of-the-art statistical analytics frameworks including Bismarck, Spark, and GraphLab. There are subtle variations to this approach: in Bismarck’s implementation, each worker processes a partition of the data, and its model is averaged at the end of each epoch; Spark implements a minibatch-based approach in which parallel workers calculate the gradient based on examples, and then gradients are aggregated by a single

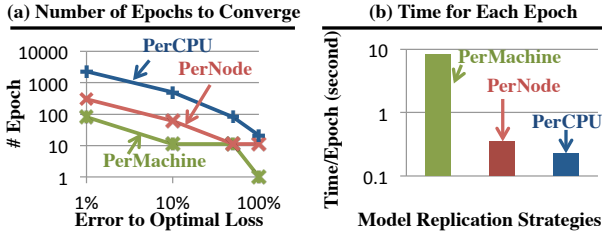


Figure 8: Illustration of Model Replication.

thread to update the final model; GraphLab implements an event-based approach where each different task is dynamically scheduled to satisfy the given consistency requirement. In DIMMWITTED, we implement PerCPU in a way that is similarly to Bismarck, where each worker has its own model replica, and each worker is responsible for updating its replica.⁴ As we will show in the experiment section, DIMMWITTED’s implementation is 3-100 \times faster than either GraphLab and Spark.⁵

PerMachine. In the PerMachine strategy, there is a single model replica that all workers update during execution. PerMachine is implemented in Hogwild! and Google’s Downpour. Hogwild! implements a lock-free protocol, which forces the hardware to deal with coherence. Although different writers may overwrite each other and readers may have dirty reads, Niu et al. [19] prove that Hogwild! converges.

PerNode. The PerNode strategy is a hybrid of PerCPU and PerMachine that takes advantage of the topological structure of NUMA machines. In PerNode, each NUMA node has a single model replica that is shared among all CPUs on that node.

Model Synchronization. Deciding how often the replicas synchronize is key to the design. In Hadoop-based and Bismarck-based models, they synchronize at the end of each epoch. This is a shared-nothing approach that works well in user-defined aggregations. However, we consider finer granularities of sharing. In DIMMWITTED, we chose to have one thread that periodically reads models on all other CPUs, averages their results, and updates each replica.

One key question for model synchronization is *how frequent should the model be synchronized?* Intuitively, we might expect that more frequent synchronization will lower the throughput; on the other hand, the more frequently we synchronize, the fewer number of iterations we might need to converge. However, in DIMMWITTED, we find that the optimal choice is to communicate as frequently as possible. The intuition is that the QPI has staggering band-

⁴We implemented MLib’s minibatch in DIMMWITTED. We find that the Hogwild!-like implementation always dominates the minibatch implementation. DIMMWITTED’s column-wise implementation for PerMachine is similar to GraphLab, with the only difference that DIMMWITTED does not schedule the task in an event-driven way.

⁵Of course, this comparison is not quite apples to apples. Note that both systems have additional sources of overhead that DIMMWITTED does not, e.g., for fault tolerance in Spark and a distributed environment in both. Both of these are unnecessary overheads in a shared-memory machine.

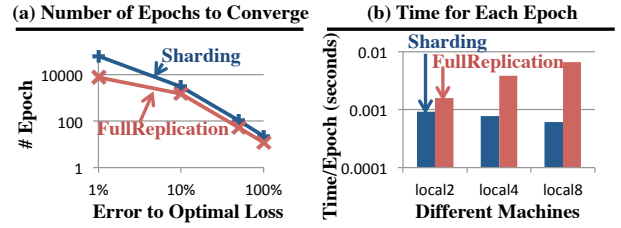


Figure 9: Illustration of Data Replication.

width (25GB/s) compared to the small amount of data we are shipping (megabytes). As a result, in DIMMWITTED, we implement an asynchronous version of the model averaging protocol: a separate thread averages models, with the effect of batching many writes together across the core into one write, reducing the number of stalls.

Tradeoffs. We observe that PerNode is more hardware efficient, as it takes less time to execute an epoch than PerMachine; PerMachine might use fewer number of epochs to converge than PerNode.

Statistical Efficiency. We observe that PerMachine usually takes fewer epochs to converge to the same loss compared to PerNode, and PerNode uses fewer number of epochs than PerCPU. To illustrate this observation, Figure 8(a) shows the number of epochs that each strategy requires to converge to a given loss for SVM (RCV1). We see that PerMachine always uses the least number of epochs to converge to a given loss: intuitively, the single model replica has more information at each step, which means there is less redundant work. We observe similar phenomena when comparing PerCPU and PerNode.

Hardware Efficiency. We observe that PerNode uses much less time to execute an epoch than PerMachine. To illustrate the difference of the time that each model replication strategy used to finish one epoch, we show in Figure 8(b) the execution time of three strategies on SVM (RCV1). We see that PerNode is 23 \times faster than PerMachine, and PerCPU is 1.5 \times faster than PerNode. PerNode takes advantage of the locality provided by the NUMA architecture. Using PMUs, we find that PerMachine incurs 11 \times more cross-node DRAM requests than PerNode.

3.4 Data Replication

In DIMMWITTED, each worker processes a subset of data and then updates its model replica. To assign a subset of data to each worker, we consider two strategies.

Sharding. Sharding is a popular strategy implemented in systems like Hogwild!, Spark, and Bismarck, in which the data set is partitioned, and each worker only works on its partition of data. When there is a single model replica, Sharding avoids wasted computation as each tuple is processed one per epoch. However, when there are multiple model replicas, Sharding might increase the variance of the estimate we form on each node, lowering the statistical efficiency. In DIMMWITTED, we implement Sharding by randomly partitioning the rows (resp. columns) of a data matrix for row-wise (resp. column-wise) access method. In column-to-row access, we also replicate other rows that are needed.

FullReplication. A simple alternative to **Sharding** is **FullReplication**, in which we replicate the whole data set many times (**PerCPU** or **PerNode**). In **PerNode**, each NUMA node will have a full copy of the data. Each node accesses its data in a different order, which means that the replicas provide non-redundant statistical information. Statistically, there are two benefits of **FullReplication**: (1) averaging different estimates from each node has a lower variance, and (2) the estimate at each node has lower variance than in the **Sharding** case, as each node’s estimate is based on the whole data. From a hardware efficiency perspective, reads are more frequent from local NUMA memory in **PerNode** than in **PerMachine**. The **PerNode** approach dominates the **PerCPU** approach as reads from the same node go to the same NUMA memory. Thus, we do not consider **PerCPU** replication from this point on.

Tradeoffs. Not surprisingly, we observe that **FullReplication** takes more time for each epoch than **Sharding**. However, we also observe that **FullReplication** uses fewer epochs than **Sharding**, especially to achieve low error. We illustrate these two observations by showing the result of running SVM on Reuters using **PerNode** in Figure 9.

Statistical Efficiency. **FullReplication** uses fewer epochs, especially to low-error tolerance. Figure 9(a) shows the number of epochs that each strategy takes to converge to a given loss. We see that for within 1% of the loss, **FullReplication** uses 10× fewer epochs on a 2-node machine. This is because each model replica sees more data than **Sharding**, and therefore has a better estimate. Because of this difference in the number of epochs, **FullReplication** is 5× faster in wall-clock time than **Sharding** to converge to 1% loss. However, we also observe that at high-error regions, **FullReplication** uses more epochs than **Sharding**, and causes a comparable execution time to a given loss.

Hardware Efficiency. Figure 9(b) shows the time for each epoch across different machines with different numbers of nodes. Because we are using the **PerNode** strategy, which is the optimal choice for this data set, the more nodes a machine has, the slower **FullReplication** is for each epoch. The slow-down is roughly consistent with the number of nodes on each machine. This is not surprising because each epoch of **FullReplication** processes more data than **Sharding**.

3.5 Implementation Details

We describe important implementation details of DIMMWITTED. In DIMMWITTED, we implement optimizations that are part of scientific computation and analytics systems. While these optimizations are not new, they are not universally implemented. However, they are critical for performance. We briefly describes each optimization and its impact.

Data and Worker Collocation. We observe that different strategies of locating data and workers affect the performance of DIMMWITTED. One standard technique is to collocate the worker and the data on the same NUMA node. In this way, the worker in each node will pull data from its own DRAM region, and does not need to occupy the node-DRAM bandwidth of other nodes. In DIMMWITTED, we tried two different placement strategies for data and workers. The first protocol, called **OS**, relies on the operating

system to allocate data and threads for workers. The operating system will usually locate data on one single NUMA node, and worker threads to different NUMA nodes using heuristics that are not exposed to the user. The second protocol, called **NUMA**, evenly distributes worker threads across NUMA nodes, and for each worker, replicates the data on the same NUMA node. We find that for SVM on RCV1, the strategy **NUMA** can be up to 2× faster than **OS**. Here are two reasons for this improvement. First, by locating data on the same NUMA node to workers, we achieve 1.24× improvement on the throughput of reading data. Second, by not asking the operating system to allocate workers, we actually have a more balanced allocation of workers on NUMA nodes.

Dense and Sparse. For statistical analytics workloads, it is not uncommon for the data matrix A to be sparse, especially for applications such as information extraction and text mining. We observe that different strategies of storing a matrix (or vector) result in different performance, depending on the sparsity of the data. In DIMMWITTED, we implement two protocols, **Dense** and **Sparse**, which store the data matrix A as a dense or sparse matrix, respectively. A **Dense** storage format has two advantages: (1) if storing a fully dense vector, it requires $\frac{1}{2}$ the space as a sparse representation, and (2) **Dense** is able to leverage hardware SIMD instructions, which allows multiple floating point operations to be performed in parallel. While a **Sparse** storage format is not able to take advantage of SIMD operations, it can save storage space and therefore improve cache performance and memory throughput. We find on a synthetic data set generated following Ji et al. [16] that when we vary the sparsity from 0.01 to 1.0, **Dense** can be up to 2× faster than **Sparse** (for sparsity=1.0) while **Sparse** can be up to 4× faster than **Dense** (for sparsity=0.01).

Row-major and Column-major Storage. There are two well-studied strategies to store a data matrix A : **Row-major** and **Column-major** storage. Not surprisingly, we observed that choosing an incorrect data storage strategy can cause a large slowdown. We conduct a simple experiment where we multiply a matrix and a vector using row-access method, where the matrix is stored in column- and row-major order. We find that the **Column-major** could 9× more L1 data load misses than using **Row-major** for two reasons: (1) our architectures fetch four doubles in a cacheline, only one of which is useful for the current operation. The prefetcher in Intel machines does not prefetch across page boundaries, and so it is unable to pick up significant portions of the strided access; (2) On the first access, the Data cache unit (DCU) prefetcher also gets the next cacheline compounding the problem, and so it runs 8× slower.⁶

4. EXPERIMENTS

We validate that exploiting the tradeoff space that we described enables DIMMWITTED’s orders of magnitude speedup over state-of-the-art competitor systems. We also validate that each tradeoff discussed in this paper affects the performance of DIMMWITTED.

⁶www.intel.com/content/dam/www/public/us/en/documents/manuals/64-ia-32-architectures-optimization-manual.pdf

Dataset		Within 1% of the Optimal Loss or RMSE				Within 50% of the Optimal Loss or RMSE			
		GraphLab	Spark	Hogwild!	DIMMWITTED	GraphLab	Spark	Hogwild!	DIMMWITTED
SVM	Reuters	> 60.0	31.4	0.1	0.1	> 60.0	7.2	0.01	0.01
	RCV1	> 300.0	> 300.0	61.4	26.8	> 300.0	> 300.0	0.71	0.17
	Music	> 300.0	> 300.0	52.2	11.2	> 300.0	> 300.0	2.48	0.27
	Forest	> 300.0	> 300.0	19.5	11.3	> 300.0	> 300.0	0.75	0.50
LR	Reuters	> 60.0	> 60.0	0.1	0.1	> 60.0	> 60.0	0.03	0.03
	RCV1	> 300.0	> 300.0	38.7	19.8	> 300.0	> 300.0	0.82	0.20
	Music	> 300.0	> 300.0	42.6	9.7	> 300.0	> 300.0	3.54	0.42
	Forest	> 300.0	> 300.0	19.6	15.6	> 300.0	> 300.0	0.68	0.52
LS	Reuters	> 60.0	> 60.0	4.1	3.2	> 60.0	57	0.17	0.09
	RCV1	> 300.0	> 300.0	27.5	10.5	> 300.0	> 300.0	1.30	0.40
	Music	> 300.0	> 300.0	8.1	4.2	> 300.0	> 300.0	0.94	0.53
	Forest	> 300.0	> 300.0	13.6	11.2	> 300.0	> 300.0	0.89	0.53
LP	Amazon	2.7	> 120.0	> 120.0	0.94	2.7	> 120.0	1.86	0.94
	Google	13.4	> 120.0	> 120.0	12.56	2.3	> 120.0	3.04	2.02
QP	Amazon	6.8	> 120.0	> 120.0	1.8	6.8	> 120.0	> 120.00	1.50
	Google	12.4	> 120.0	> 120.0	4.3	9.9	> 120.0	> 120.00	3.70

Figure 11: End-to-End Comparison (time in seconds). We take 5 runs on local2 and report the average (standard deviation for all numbers < 5% of the mean). Entries with > indicate a timeout.

Model	Dataset	#Row	#Col.	NNZ	Size (Sparse)	Size (Dense)	Sparse
SVM LR LS	RCV1	781K	47K	60M	914MB	275GB	✓
	Reuters	8K	18K	93K	1.4MB	1.2GB	✓
	Music	515K	91	46M	701MB	0.4GB	
	Forest	581K	54	30M	490MB	0.2GB	
LP	Amazon	926K	335K	2M	28MB	>1TB	✓
	Google	2M	2M	3M	25MB	>1TB	✓
QP	Amazon	1M	1M	7M	104MB	>1TB	✓
	Google	2M	2M	10M	152MB	>1TB	✓
Gibbs	Paleo	69M	30M	108M	2GB	>1TB	✓
NN	MNIST	120M	800K	120M	2GB	>1TB	✓

Figure 10: Dataset Statistics. NNZ refers to the Number of Non-zero elements. The # columns are also equal to the number of variables in the model.

4.1 Experiment Setup

We describe the details of our experimental setting.

Datasets and Statistical Models. We validate the performance and quality of DIMMWITTED on a diverse set of statistical models and datasets. For statistical models, we choose five models that are among the most popular models used in statistical analytics: (1) Support Vector Machine (SVM), (2) Logistic Regression (LR), (3) Least Squares Regression (LS), (4) Linear Programming (LP), and (5) Quadratic Programming (QP). For each model, we choose data sets with different characteristics, including size, sparsity, and under- or over-determination. For SVM, LR, and LS, we choose four data sets: Reuters⁷, RCV1⁸, Music⁹, and Forest.¹⁰ Reuters and RCV1 are data sets for text classification that are sparse and underdetermined. Music and Forest are standard benchmark datasets that are dense and over-determined. For QP and LR, we consider a social-network ap-

plication, i.e., network analysis, and use two data sets from Amazon’s customer data and Google’s Google+ social networks.¹¹ Figure 10 shows the data set statistics.

Metrics. We measure the quality and performance of DIMMWITTED and other competitors. To measure the quality, we follow prior art and use RMSE for SVM, and the loss function for LR, LS, QP and LP. For end-to-end performance, we measure the wall-clock time it takes for each system to converge to a loss (or RMSE) that is within 100%, 50%, 10%, and 1% of the optimal loss (or RMSE).¹² When measuring the wall-clock time, we do not count the time used for data loading and result outputting for all systems. We also use other measurements to understand the details of the tradeoff space, including (1) Local LLC request, (2) Remote LLC request, and (3) Local DRAM request. We use Intel Performance Monitoring Units (PMUs) and follow its manual¹³ to conduct these experiments.

Experiment Setting. We compare DIMMWITTED with three competitor systems: GraphLab [17], Spark [29], and Hogwild! [19]. GraphLab is a distributed graph processing system that supports a large range of statistical models. Spark is an in-memory implementation of the MapReduce framework that provides a package called MLlib for machine learning algorithms. Hogwild! is an in-memory lock-free framework for statistical analytics. We find that all three systems pick some points in the tradeoff space that we considered in DIMMWITTED. In GraphLab, all models are implemented using stochastic coordinate descent (column-wise access); in Spark/MLlib and Hogwild!, SVM and LR are implemented using stochastic gradient descent (row-wise access). We use implementations that are provided by the original developers whenever possible.

We run experiments on a variety of representative architectures. These machines differ in a range of configurations, including number of NUMA nodes, the size of last-level

⁷archive.ics.uci.edu/ml/datasets/Reuters-21578+Text+Category+Collection

⁸about.reuters.com/researchandstandards/corpus/

⁹archive.ics.uci.edu/ml/datasets/YearPredictionMSD

¹⁰archive.ics.uci.edu/ml/datasets/Coverttype

¹¹snap.stanford.edu/data/

¹²We get the optimal loss (or RMSE) by running all systems for one hour and choose the lowest.

¹³software.intel.com/en-us/articles/performance-monitoring-unit-guidelines

cache (LLC), and memory bandwidth. See Figure 3 for a summary of these machines. DIMMWITTED, Hogwild!, and GraphLab are implemented using C++ and Spark is implemented using Scala. We tune both GraphLab and Spark according to their best practice guidelines.¹⁴ Systems are compiled with g++ 4.7.2 (-O3), Java 1.7, or Scala 2.9.

4.2 End-to-End Comparison

We validate that DIMMWITTED outperforms competitor systems in terms of end-to-end performance and quality.

Protocol. For each system, we grid search their statistical parameters, including step size ($\{0.1, 0.01, 0.001, 0.0001\}$) and mini-batch size for Spark ($\{\{1\%, 10\%, 50\%, 100\%\}\}$); we always report the best configuration, which is essentially the same for each system. We measure the time it takes for each system to find a solution that is within 1%, 10%, and 50% of the optimal loss or RMSE. Figure 11 shows the result for 1% and 50%; the results for 10% are similar. We report end-to-end numbers from local2 that has 2 nodes and 24 logical nodes, as GraphLab does not run on machines with more than 64 logical cores. Figure 14 shows the DIMMWITTED’s choice of point in the tradeoff space on local2.

As shown in Figure 11, DIMMWITTED always converges to the given loss in less time than the other competitors. On SVM and LR, DIMMWITTED could be up to 10 \times faster than Hogwild!, and more than two orders of magnitude faster than GraphLab and Spark. The difference between DIMMWITTED and Hogwild! is greater for LP and QP, where DIMMWITTED outperforms Hogwild! by more than two orders of magnitude. On LP and QP, DIMMWITTED is also up to 3 \times faster than GraphLab, and two orders of magnitude faster than Spark.

Tradeoff Choices. We dive more deeply into these numbers to substantiate our claim that there are some points in the tradeoff space that are not used by GraphLab, Hogwild!, and Spark. Each tradeoff selected by our system is in Figure 14. For example, GraphLab uses column-wise access for all models, while Spark and Hogwild! use row-wise access for all models and allow only PerMachine model replication. These special points work well for some—but not all—models. For example, for LP and QP, GraphLab is only 3 \times slower than DIMMWITTED, which chooses column-wise and PerMachine. This factor of 3 is to be expected as GraphLab also allows distributed access and so has additional overhead. But there are other points: on SVM and LR, DIMMWITTED outperforms GraphLab, because the column-wise algorithm implemented by GraphLab is not as efficient as row-wise on the same data set. DIMMWITTED outperforms Hogwild! because DIMMWITTED takes advantage of model replication, while Hogwild! incurs 11 \times more cross-node DRAM requests than DIMMWITTED; in contrast, DIMMWITTED incurs 11 \times more local DRAM requests than Hogwild! does.

DIMMWITTED outperforms Spark partly due to algorithmic reasons: Spark parallelizes SGD using mini-batches and incurs significant overhead to synchronize after each mini-batch; thus DIMMWITTED finishes each epoch 10 times faster. However, on LP and QP, DIMMWITTED outperforms Spark

	SVM (RCV1)	LR (RCV1)	LS (RCV1)	LP (Google)	QP (Google)	Parallel Sum
GraphLab	0.2	0.2	0.2	0.2	0.1	0.9
Spark	0.1	0.1	0.1	0.05	0.01	0.2
Hogwild!	1.3	1.4	1.3	0.3	0.2	13
DIMMWITTED	5.1	5.2	5.2	0.7	1.3	21

Figure 13: Comparison of Throughput (GB/seconds) of Different Systems on local2.

		Access Methods	Model Replication	Data Replication
SVM	Reuters	Row-wise	PerNode	FullReplication
LR	RCV1			
LS	Music			
LP	Amazon	Column-wise	PerMachine	FullReplication
QP	Google			

Figure 14: Plans that DimmWitted Chooses in the Tradeoff Space for Each Data Set on Machine local2.

and Hogwild! because the row-wise access method implemented by these systems are not as efficient as column-wise access on the same data set. DIMMWITTED outperforms GraphLab because DIMMWITTED finishes each epoch up to 3 \times faster.

Throughput. We compare the throughput of different systems for an extremely simple task: parallel sums. Our implementation of parallel sum follows our implementation of other statistical models (with a trivial update function), and uses all cores on a single machine. Figure 13 shows the throughput on all systems on different models on one data set. We see from Figure 13 that DIMMWITTED achieves the highest throughput of all the systems. For parallel sum, DIMMWITTED is 1.6 \times faster than Hogwild!, and we find that DIMMWITTED incurs 8 \times fewer LLC cache misses than Hogwild!. Compared with Hogwild!, in which all threads write to a single copy of the sum result, DIMMWITTED maintains one single copy of the sum result per NUMA node, and therefore, the workers on one NUMA node do not invalidate the cache on another NUMA node. When running on only a single thread, DIMMWITTED has the same implementation as Hogwild!. Compared with GraphLab, DIMMWITTED is 20 \times faster, likely due to the overhead of GraphLab dynamically scheduling tasks. To compare with Spark, which is written in Scala, we implemented a Scala version, which is 3 \times slower than C++; this suggests that the overhead is not just due to the language, and the overhead of scheduling is a factor in Spark as well.

4.3 Tradeoffs of DIMMWITTED

We validate that all tradeoffs we describe in this paper have an impact on the efficiency of DIMMWITTED. We report on a more modern architecture, local4 with 4 NUMA sockets, in this section. We describe how the results change with different architecture.

4.3.1 Access Method Selection

We validate that different access methods have different performance, and that no single access method dominates

¹⁴Spark: spark.incubator.apache.org/docs/0.6.0/tuning.html; GraphLab: graphlab.org/tutorials-2/fine-tuning-graphlab-performance/.

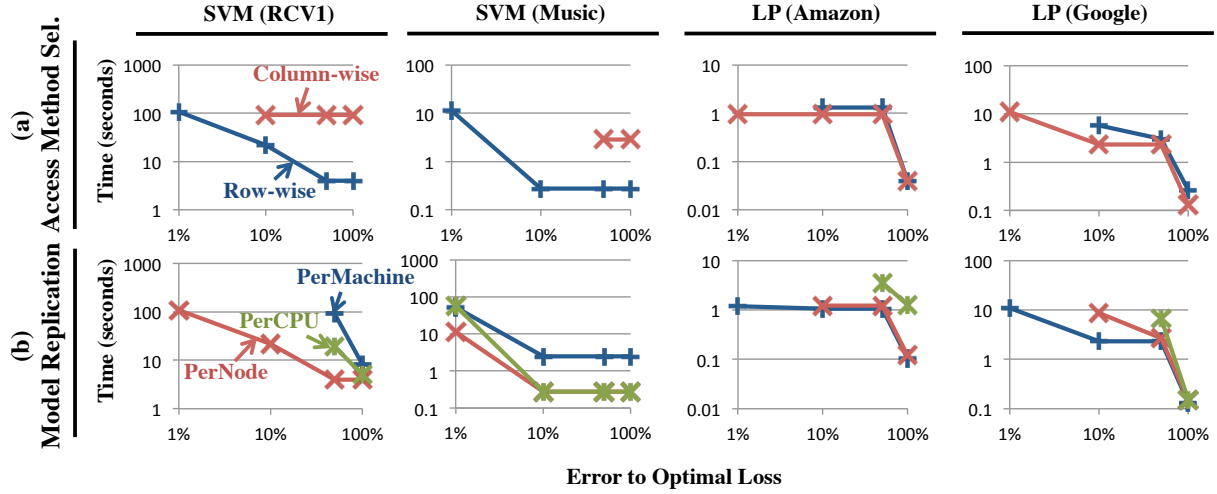


Figure 12: Tradeoffs in DimmWitted. Missing points timeout in 120 seconds.

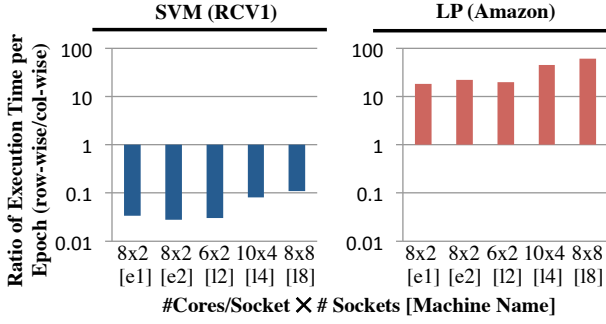


Figure 15: Ratio of Execution Time per Epoch (row-wise/column-wise) on Different Architectures. A number larger than 1 means row-wise is slower. 12 means local2, e1 means ec2.1, etc.

the others. We run DIMMWITTED on all statistical models and compare two strategies, row-wise and column-wise. In each experiment, we force DIMMWITTED to use the corresponding access method, but report the best point for the other tradeoffs. Figure 12(a) shows the result as we measure the time it takes to achieve each loss. The more stringent loss requirements (1%) are on the left-hand side. The horizontal line segments in the graph indicate that a model may reach, say, 50% as quickly (in epochs) as it reaches 100% loss.

We see from Figure 12(a) that the difference between row-wise and column-to-row access could be more than 100 \times for different models. For SVM on RCV1, row-wise access converges at least 4 \times faster to 10% loss and at least 10 \times faster to 100% loss. We observe similar phenomena on Music; compared with RCV1, column-to-row access converges to 50% loss and 100% loss at a 10 \times slower rate. On such data sets, the column-to-row access requires simply more reads and writes. This supports the folk wisdom that gradient methods are preferable to coordinate descent methods. On the other hand, for LP, column-wise access dominates: row-wise access does not converge to 1% loss within the timeout period for either Amazon or Google. Column-wise access converges at least 10-100 \times faster than row-wise access to 1%

loss. We also observe that LR is similar to SVM, and QP is similar to LP. Thus, no single access method dominates the other.

The cost of writing and reading are different, and is captured by a parameter that we called α in Section 3.2. We describe the impact of this factor on the relative performance of row- and column-wise strategies. Figure 15 shows the ratio of the time that each strategy used (row-wise/column-wise) on SVM(RCV1) and LP(Amazon). We see that, as the number of sockets on a machine increases, the ratio of execution time gets larger, which means that row-wise gets slower relative to column-wise, i.e., with increasing α . As the write cost captures the cost of a hardware-resolved conflict, we see that this constant is likely to grow. Thus, if next generation architectures increase in number of sockets, the cost parameter α and consequently the importance of this tradeoff are likely to grow.

Cost-based Optimizer. We observed that on all data sets, our cost-based optimizer selects row-wise access for SVM, LR, and LS, and column-wise access for LP and QP. These choices are consistent with what we observed in Figure 12.

4.3.2 Model Replication

We validate that there is no single strategy for model replication that dominates the others. We force DIMMWITTED to run strategies in PerMachine, PerNode, and PerCPU and choose other tradeoffs by choosing the plan that achieves the best result. Figure 12(b) shows the result.

We see from Figure 12(b) that the gap between PerMachine and PerNode could be up to 100 \times . We first observe that PerNode dominates PerCPU on all data sets. For SVM on RCV1, PerNode converges 10 \times faster than PerCPU to 50% loss, and on other models and data sets, we observe a similar phenomenon. This is due to the low statistical efficiency of PerCPU, as we discussed in Section 3.3. Although PerCPU eliminates write contention inside one NUMA node, this write contention is less critical. On large models and machines with small caches, we have also observed that PerCPU could spill the cache.

These graphs show that neither PerMachine nor PerNode dominates the other across all data sets and statistical mod-

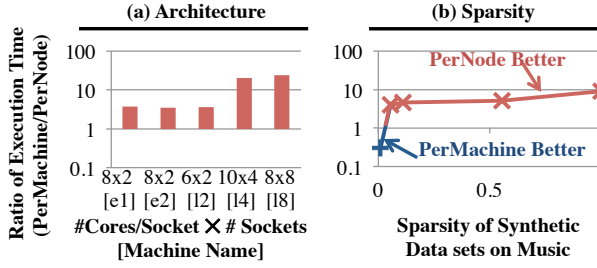


Figure 16: The Impact of Different Architectures and Sparsities on Model Replication. A ratio larger than 1 means that **PerNode** converges faster than **PerMachine** to 50% loss.

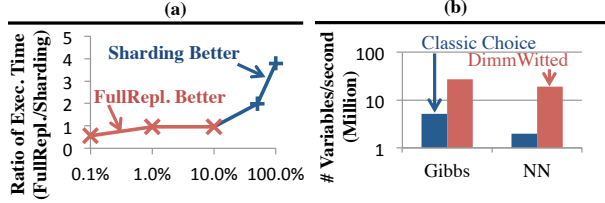


Figure 17: (a) Tradeoffs of Data Replication. A ratio smaller than 1 means that **FullReplication** is faster. (b) Performance of Gibbs Sampling and Neural Networks Implemented in **DimmWitted**.

els. For SVM on RCV1, **PerNode** converges $12\times$ faster than **PerMachine** to 50% loss. However, for LP on Amazon, **PerMachine** is at least $14\times$ faster than **PerNode** to converge to 1% loss. For SVM, the reason that **PerNode** converges faster is because it has $5\times$ higher throughput than **PerMachine**, and for LP, the reason that **PerNode** is slower is because **PerMachine** takes at least $10\times$ fewer epochs to converge to a small loss. One interesting observation is that for LP on Amazon, **PerMachine** and **PerNode** do have comparable performance to converge to 10% loss. Compared with the 1% loss case, this implies that **PerNode**'s statistical efficiency decreases as the algorithm tries to achieve a smaller loss. This is not surprising, as more effort is put into reconciling the **PerNode** estimates.

We observe that the relative performance of **PerMachine** and **PerNode** depends on (1) the number of sockets used on each machine, and (2) the sparsity of the update.

To validate (1), we measure the time that **PerNode** and **PerMachine** take on SVM(RCV1) to converge to 50% loss on various architectures, and we report the ratio (**PerMachine**/**PerNode**) in Figure 16. We see that **PerNode**'s relative performance improves with the number of sockets. We attribute this to the increased cost of write contention in **PerMachine**.

To validate (2), we generate a series of synthetic data sets each of which subsamples the elements in each row of the Music dataset; Figure 16(b) shows the result. When the sparsity is 1%, **PerMachine** outperforms **PerNode**, as each update touches only one element of the model; thus, the write contention in **PerMachine** is not a bottleneck. As the sparsity increases (i.e., the update gets more dense), we observe that **PerNode** outperforms **PerMachine**.

4.3.3 Data Replication

We validate the impact of different data replication strategies. We run **DIMMWITTED** by fixing data replication strategies to **FullReplication** or **Sharding**, and choose the best plan for each other tradeoffs. We measure the execution time for each strategy to converge to a given error for SVM on the same data set, RCV1. We report the ratio of these two strategies as (**FullReplication**/**Sharding**) in Figure 17(a). We see that for low-error region (e.g., 0.1%), **FullReplication** is $1.8\text{--}2.5\times$ faster than **Sharding**. This is because **FullReplication** decreases the skew of data assignment to each worker, and hence each individual model replica can form a more accurate estimate. For the high-error region (e.g., 100%), we observe that **FullReplication** appears to be $2\text{--}5\times$ slower than **Sharding**. We find that, for 100% loss, both **FullReplication** and **Sharding** converge in a single epoch and **Sharding** may therefore be preferred, as it examines less data to complete that single epoch. In all of our experiments, **FullReplication** is never substantially worse and can be dramatically better. Thus, if there is available memory, the **FullReplication** data replication seems to be clearly preferable.

5. EXTENSIONS

We describe how to run Gibbs sampling (which uses a column-to-row access method) and deep neural networks (which use a row access method). Using the same tradeoffs (and engine), we achieve significant increase in speed over classical implementation choices of these algorithms.

5.1 Gibbs Sampling

Gibbs sampling is one of the most popular algorithms to solve statistical inference and learning over probabilistic graphical models [23]. We briefly describe Gibbs sampling over factor graph and observe its main step is a column-to-row access. A factor graph can be thought of as a bipartite graph of a set of variables and a set of factors. To run Gibbs sampling, the main operation is to select a single variable, and calculate the conditional probability of this variable, which requires fetching all factors that contain this variable and all assignments of variables connected to these factors. This operation corresponds to the column-to-row access method. Similar to first order methods, recently a Hogwild! algorithm for Gibbs was established [11]. As shown in Figure 17(b), applying the technique in **DIMMWITTED** to Gibbs sampling has $4\times$ the throughput of samples as the **PerMachine** strategy.

5.2 Deep Neural Networks

Neural networks are one of the most classic machine learning models [18]; recently these models have been intensively revisited by adding more layers [6, 13]. A deep neural network contains multiple layers in which each layer contains a set of neurons (variables). Different neurons connect with each other only by links across consecutive layers. The value of one neuron is a function of all other neurons in the previous layer and a set of weights. Variables in the last layer have human labels as training data; the goal of deep neural network learning is to find the set of weights that maximizes the likelihood of the human labels. Back-propagation with stochastic gradient descent is the de facto method of optimizing a deep neural network.

Following LeCun et al. [14], we implement SGD over a 7-layer neural network with 0.12 billion neurons and 0.8 million parameters, using a standard handwriting-recognition

benchmark data set called MNIST¹⁵. Figure 17(b) shows the number of variables (neurons) that are processed by DIMMWITTED per second. For this application, DIMMWITTED uses *PerNode* and *FullReplication*, and the classical choice made by LeCun is *PerMachine* and *Sharding*. As shown in Figure 17(b), DIMMWITTED achieves more than an order of magnitude higher throughput than this classic baseline (to achieve the same quality as reported in this classical paper).

6. RELATED WORK

We review work in three main areas: statistical analytics, main-memory databases, and mathematical optimization.

Statistical Analytics. There is a trend to integrate statistical analytics into data processing systems. Database vendors have recently put our new products in this space, including Oracle, Pivotal’s MADlib [8], and IBM’s SystemML [7], and SAP’s HANA. These systems support statistical analytics in existing data management systems. A key challenge for statistical analytics is performance.

Most data processing frameworks developed in the last few years are designed to support statistical analytics including Mahout for Hadoop, MLI for Spark [26], GraphLab [17], and MADlib for PostgreSQL or Greenplum [8]. Although these systems increase the performance of corresponding statistical analytics tasks significantly, we observed that each of them only implements some points in DIMMWITTED’s trade-off space. Compared with these systems, DIMMWITTED studies a larger tradeoff space and designs guidelines and optimizers to select among different strategies.

Main-memory Databases. The database community has recognized that multi-socket, large-memory machines have changed the data processing landscape, and there is a flurry of recent work about how to build in-memory analytics systems [2, 3, 5, 12, 15, 20, 21, 28]. Classical tradeoffs have been revisited on the modern architecture to gain significant improvement: Balkesen et al. [3], Albutiu et al. [2], Kim et al. [12], and Li [15] studied the tradeoff for joins and shuffling respectively. This work takes advantage of modern architectures, e.g., NUMA and SIMD, to increase memory bandwidth; we study a new tradeoff space for statistical analytics in which the performance of the system is affected by both hardware efficiency and statistical efficiency.

Mathematical Optimization. Many statistical analytics tasks are mathematical optimization problems. Recently, the mathematical optimization community has been looking at how to parallelize optimization problems [16, 19, 30]. For example, Niu et al. [19] and Ji et al. [16] develop the lock-free protocol for either SGD and SCD, respectively.

7. CONCLUSION

For statistical analytics on main-memory, NUMA-aware machines, we studied tradeoffs in access methods, model replication, and data replication. We found that using novel points in this tradeoff space can have a substantial benefit: our DIMMWITTED prototype engine can run at least one popular task at least 100× faster than other competitor systems. This comparison demonstrates that this tradeoff

space may be interesting for the current and next generation of statistical analytics systems.

Acknowledgements We would like to thank Arun Kumar, Victor Bittorf, the teams at Advanced Analytics at Oracle, and Greenplum/Pivotal, and Impala’s Cloudera team for their support and for sharing their experiences in building analytics systems.

8. REFERENCES

- [1] A. Agarwal, O. Chapelle, M. Dudík, and J. Langford. A reliable effective terascale linear learning system. *ArXiv e-prints*, 2011.
- [2] M.-C. Albutiu, A. Kemper, and T. Neumann. Massively parallel sort-merge joins in main memory multi-core database systems. *PVLDB*, pages 1064–1075, 2012.
- [3] C. Balkesen and et al. Multi-core, main-memory joins: Sort vs. hash revisited. *PVLDB*, pages 85–96, 2013.
- [4] L. Bergstrom. Measuring NUMA effects with the STREAM benchmark. *ArXiv e-prints*, 2011.
- [5] C. Chasseur and J. M. Patel. Design and evaluation of storage organizations for read-optimized main memory databases. *PVLDB*, pages 1474–1485, 2013.
- [6] J. Dean and et al. Large scale distributed deep networks. In *NIPS*, pages 1232–1240, 2012.
- [7] A. Ghoting and et al. SystemML: Declarative machine learning on MapReduce. In *ICDE*, pages 231–242, 2011.
- [8] J. M. Hellerstein and et al. The MADlib analytics library: Or MAD skills, the SQL. *PVLDB*, pages 1700–1711, 2012.
- [9] M. Hornick. Low-rank matrix factorization in Oracle R Advanced Analytics for Hadoop. https://blogs.oracle.com/R/entry/low_rank_matrix_factorization_in.
- [10] Impala. <http://www.cloudera.com/content/cloudera/en/products-and-services/cdh/impala.html>.
- [11] M. J. Johnson, J. Saunderson, and A. S. Willsky. Analyzing Hogwild parallel Gaussian Gibbs sampling. In *NIPS*, pages 2715–2723, 2013.
- [12] C. Kim and et al. Sort vs. hash revisited: Fast join implementation on modern multi-core CPUs. *PVLDB*, pages 1378–1389, 2009.
- [13] Q. V. Le and et al. Building high-level features using large scale unsupervised learning. In *ICML*, pages 8595–8598, 2012.
- [14] Y. LeCun and et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, pages 2278–2324, 1998.
- [15] Y. Li and et al. NUMA-aware algorithms: the case of data shuffling. In *CIDR*, 2013.
- [16] J. Liu and et al. An asynchronous parallel stochastic coordinate descent algorithm. *ArXiv e-prints*, 2013.
- [17] Y. Low and et al. Distributed GraphLab: A framework for machine learning in the cloud. *PVLDB*, pages 716–727, 2012.
- [18] T. M. Mitchell. *Machine Learning*. McGraw-Hill, USA, 1997.
- [19] F. Niu and et al. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In *NIPS*, pages 693–701, 2011.
- [20] L. Qiao and et al. Main-memory scan sharing for multi-core CPUs. *PVLDB*, pages 610–621, 2008.
- [21] V. Raman and et al. DB2 with BLU acceleration: So much more than just a column store. *PVLDB*, pages 1080–1091, 2013.
- [22] P. Richtárik and M. Takáč. Parallel coordinate descent methods for big data optimization. *ArXiv e-prints*, 2012.
- [23] C. P. Robert and G. Casella. *Monte Carlo Statistical Methods (Springer Texts in Statistics)*. Springer, USA, 2005.
- [24] A. Smola and S. Narayanamurthy. An architecture for parallel topic models. *PVLDB*, pages 703–710, 2010.
- [25] S. Sonnenburg and et al. The SHOGUN machine learning toolbox. *J. Mach. Learn. Res.*, pages 1799–1802, 2010.
- [26] E. Sparks and et al. MLI: An API for distributed machine learning. In *ICDM*, pages 1187–1192, 2013.
- [27] S. Sridhar and et al. An approximate, efficient LP solver for LP rounding. In *NIPS*, pages 2895–2903, 2013.
- [28] S. Tu and et al. Speedy transactions in multicore in-memory databases. In *SOSP*, pages 18–32, 2013.
- [29] M. Zaharia and et al. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *NSDI*, 2012.
- [30] M. Zinkevich and et al. Parallelized stochastic gradient descent. In *NIPS*, pages 2595–2603, 2010.

¹⁵yann.lecun.com/exdb/mnist/